

Phylogenetic incongruence through the lens of Monadic Second Order logic

Steven Kelk, Leo van Iersel and Celine Scornavacca

¹ Department of Knowledge Engineering (DKE), Maastricht University, P.O. Box 616, 6200 MD Maastricht, The Netherlands, steven.kelk@maastrichtuniversity.nl

² Delft Institute of Applied Mathematics, Delft University of Technology, P.O. Box 5, 2600 AA Delft, The Netherlands, l.j.j.v.iersel@gmail.com

³ Institut des Sciences de l'Evolution (Université de Montpellier, CNRS, IRD, EPHE), Place E. Bataillon - CC 064 - 34095 Montpellier Cedex 5, France, celine.scornavacca@univ-montp2.fr

Abstract. Within the field of phylogenetics there is growing interest in measures for summarising the dissimilarity, or *incongruence*, of two or more phylogenetic trees. Many of these measures are NP-hard to compute and this has stimulated a considerable volume of research into fixed parameter tractable algorithms. In this article we use *Monadic Second Order* logic (MSOL) to give alternative, compact proofs of fixed parameter tractability for several well-known incongruency measures. In doing so we wish to demonstrate the considerable potential of MSOL - machinery still largely unknown outside the algorithmic graph theory community - within phylogenetics. A crucial component of this work is the observation that many of these measures, when bounded, imply the existence of an *agreement forest* of bounded size, which in turn implies that an auxiliary graph structure, the *display graph*, has bounded treewidth. It is this bound on treewidth that makes the machinery of MSOL available for proving fixed parameter tractability. We give a variety of different MSOL formulations. Some are based on explicitly encoding agreement forests, while some only use them implicitly to generate the treewidth bound. Our formulations introduce a number of “phylogenetics MSOL primitives” which will hopefully be of use to other researchers.

1 Introduction

The central goal of phylogenetics is to accurately infer the evolutionary history of a set of species (or *taxa*) X from incomplete information. Classically, phylogenetic reconstruction has access to information about each element in X , such as DNA data, and seeks to infer a phylogenetic tree - a tree whose leaves are bijectively labeled by X - that best fits this data. There is a vast literature available on this topic and many different algorithms exist for constructing phylogenetic trees [16,26]. In practice, it is not uncommon for phylogenetic analysis to generate multiple phylogenetic trees as output. This can occur for various reasons, ranging from software engineering choices (many tree-building packages are designed to generate multiple optimal and near-optimal solutions) to more structural explanations (*reticulate* evolutionary signals that are comprised of multiple distinct tree signals). Given two (or more) distinct phylogenetic trees, it is natural to compare them to determine whether the difference is significant. This explains the interest of the phylogenetics community for measures that can quantify the dissimilarity, or *incongruence*, of phylogenetic trees [20]. Some of these measures (such as TREE BISECTION AND RECONNECTION distance [1]) are studied to better understand how local-search heuristics, based on rearrangement operations, navigate the space of phylogenetic trees (e.g., [9]). Others, such as HYBRIDIZATION NUMBER [8], are studied because they assist with the inference of phylogenetic networks, which generalise phylogenetic trees to directed acyclic graphs [20,21].

Unfortunately, many of these measures are NP-hard and APX-hard to compute. On the positive side, however, the phylogenetics community has been quite successful in proving that these measures are fixed parameter tractable (FPT) in their natural parameterizations. Informally this means that a measure that evaluates to k can be computed in time $f(k) \cdot \text{poly}(n)$ where f is some function that only depends on k and n is the size of the instance (often taken to be $|X|$). Such

running times have the potential to be much faster than running times of the form $O(n^{f(k)})$ when the measure in question is comparatively small. See e.g. [15] for more background on FPT. A number of state-of-the-art phylogenetics software packages are based on FPT algorithms, such as the software used in [28]. Most FPT results in the phylogenetics literature are based on classical proof techniques such as polynomial-time kernelization and bounded-search.

Parallel to all of this, algorithmic graph theorists have made great steps forward in identifying sufficient, structural conditions under which NP-hard problems on graphs become (fixed parameter) tractable. At the heart of this research lies the width parameter, the most famous example being *treewidth*. Informally treewidth is a measure that quantifies the dissimilarity of a graph from being a tree. The notion of treewidth, which is most famously associated with the celebrated Graph Minors project of Robertson and Seymour [24], has had a profound impact upon algorithm design. A great many NP-hard problems turn out to become tractable on graphs of bounded treewidth, using broadly similar proof techniques i.e. dynamic programming on tree decompositions [4]. This contributed to the rise of meta-theorems, the archetypal example being *Courcelle’s Theorem* [14,2]. This states, when combined with the result from [5], that any graph property that can be abstractly formulated as a length ℓ sentence of *Monadic Second Order* logic (MSOL), can be tested in time $f(t, \ell) \cdot O(n)$ on graphs of treewidth t , where n is the number of vertices in the graph. When t and ℓ are both bounded by a function of a single parameter p , this yields a running time of the form $f(p) \cdot O(n)$ i.e. linear-time fixed parameter tractability in parameter p . This is an extremely powerful technique in the sense that it completely abstracts away from ad-hoc algorithm design and permits highly compact, “declarative” proofs that a problem is FPT. Courcelle’s Theorem (and its variants) are more than two decades old, but their potential is rarely exploited by the phylogenetics community. One exception is the literature on *unrooted compatibility*, which asks whether a set of unrooted phylogenetic trees are compatible. The FPT proof by Bryant and Lagergren [10] proves that the *display graph* (the graph obtained by identifying all taxa with the same label) has bounded treewidth (in the number of input trees), and then gives an MSOL formulation which tests compatibility. A follow-up result by the present authors applies a similar approach [25].

In this article we show that this technique has much broader potential within phylogenetics. To clarify the exposition we focus on binary trees (both rooted and unrooted) on the same set of taxa X . We begin by proving that if two trees have an *agreement forest* of size k – essentially a partition of the trees into k isomorphic subtrees – the treewidth of the display graph is bounded by a function of k . This simple observation is significant because of the prominent role of agreement forests within the phylogenetics literature. We use this insight to re-analyse three well-known NP-hard phylogenetics problems that were previously shown to be FPT using more conventional analysis. In particular, we give MSOL formulations for (1) UNROOTED MAXIMUM AGREEMENT FOREST (uMAF), which is equivalent to the problem of computing TREE BISECTION AND RECONNECTION distance (TBR) on unrooted trees, (2) ROOTED MAXIMUM AGREEMENT FOREST (rMAF), which is equivalent to the problem of computing ROOTED SUBTREE PRUNE AND REGRAFT distance (rSPR) on rooted trees, and (3) HYBRIDIZATION NUMBER (HN) on rooted trees. The formulations for uMAF and rMAF are based on explicitly modelling agreement forests using quartets and edge cuts. The formulation for HN uses agreement forests implicitly to obtain the treewidth bound but, due to the difficulties in encoding *acyclic* agreement forests, then bypasses the agreement forest abstraction. Instead, it encodes an equivalent, “elimination ordering” formulation of HN which considers sequences of pruned common subtrees. Finally we consider the (4) MAXIMUM PARSIMONY DISTANCE ON BINARY CHARACTERS problem. This asks for a binary character f on X that maximizes the absolute difference between the parsimony score of f on the two trees. It is NP-hard but not known to be FPT (in the parsimony distance). Here we give an optimization MSOL formulation which shows that the problem is FPT in parameter uMAF. Although this does not settle whether the natural parameterization of the problem is FPT, it does demonstrate a number of interesting principles. Firstly, it demonstrates the power of “simulating” the execution of polynomial-time algorithms (in this case, Fitch’s algorithm [18]) within MSOL. Secondly, any subsequent proof that TBR distance is at most a bounded distance above d_{MP}^2 distance and/or

that d_{MP}^2 distance induces bounded treewidth display graphs, will automatically prove that d_{MP}^2 distance is FPT in its natural parameterization.

Summarizing, our formulations show the potential for MSOL to generate compact, logical FPT proofs for phylogenetics problems. The machinery of MSOL does not yield practical algorithms but it is an excellent classification tool. Once the existence of FPT algorithms has been confirmed via MSOL one can then switch efforts to finding a *good* FPT algorithm by more direct analysis, possibly (but not exclusively) through direct analysis of tree decompositions. Our formulations also introduce a number of phylogenetics “primitives” concerning quartets, clusters, subtrees and compatibility that we hope will be of use to other phylogenetics researchers.

2 Preliminaries

In this section, we define the main objects that will be manipulated in this paper.

An *unrooted phylogenetic tree* T (unrooted tree for short) is a tree in which no vertex has degree 2 and in which the leaves are bijectively labeled by a label set $\mathcal{L}(T)$. The leaf labels are often called *taxa* and the symbol X is frequently used as shorthand for $\mathcal{L}(T)$. Internal vertices are not labeled. A *rooted phylogenetic tree* (rooted tree for short) is defined similarly, except that it has exactly one vertex, called the *root* of the tree, that is permitted to have degree 2, and edges are directed away from the root. An unrooted tree is *binary* if every internal vertex has degree 3, and a rooted tree is binary if each internal vertex has indegree 1 and outdegree 2, and the root has outdegree 2 and indegree 0.

Given an unrooted tree T and a subset $Y \subseteq \mathcal{L}(T)$, we use $T(Y)$ to denote the minimal subtree of T connecting Y . Moreover, we denote by $T|_Y$ the tree obtained from $T(Y)$ when suppressing vertices of degree 2. We say that $T|_Y$ is the subtree of T *induced* by Y . In graph theory terms, $T|_Y$ is a label-preserving topological minor of T . Induced subtrees are defined in the same way for rooted trees, except that the root of $T|_Y$ becomes the vertex in the minimal connecting subgraph that is closest to the root of T , and we suppress all degree 2 vertices except the new root. We write $T - Y$ to denote $T|_{\mathcal{L}(T)-Y}$. For any node u of a rooted tree T , T_u is the subtree of T rooted at u .

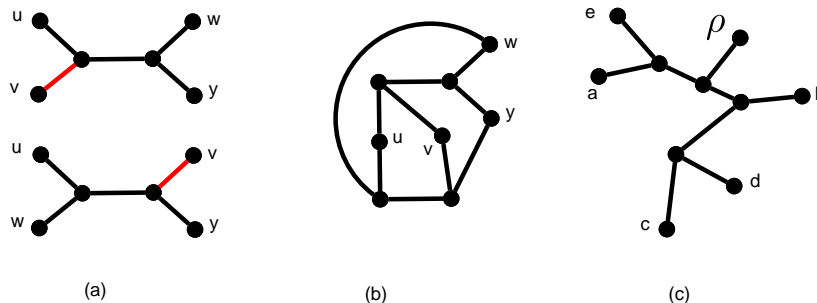


Fig. 1. (a) Two unrooted binary phylogenetic trees on $\{u, v, w, y\}$. A maximum agreement forest (uMAF) for these two trees contains 2 components, and can be obtained by cutting the single red edge in both trees and then suppressing the resulting degree 2 vertices. (b) The display graph for the two trees from (a), obtained by identifying leaves with the same label. (c) How addition of a special label ρ can be used to root a tree: edges are assumed to be directed away from ρ .

Given a label set X , a *bipartition* (or *split*) $A|B$ on X is a partition of X into two non-empty sets. Each edge $\{u, v\}$ of a tree T induces a split $\mathcal{L}(T_u)|\mathcal{L}(T_v)$, where T_u and T_v are the two trees obtained from T when $\{u, v\}$ is deleted. Given a rooted tree T with label set X , a subset X' of X is called a *clade* (or *cluster*) of T , if T contains a node v such that $\mathcal{L}(T_v) = X'$.

Given an unrooted binary tree T and a set of four distinct labels $\{u, v, w, y\}$ in $\mathcal{L}(T)$, $T|_{\{u, v, w, y\}}$ will be exactly one of the three possible unrooted binary trees on $\{u, v, w, y\}$. These are called

quartets and are denoted respectively by $uv|wy$, $uw|vy$ and $wv|uy$, depending on the bipartition induced by its central edge. In Figure 1(a) we see $uv|wy$ and $uw|vy$. Given a rooted binary tree T and a set of three labels $\{u, v, w\}$ in $\mathcal{L}(T)$, $T|_{\{u,v,w\}}$ will be exactly one of the three possible rooted binary trees on $\{u, v, w\}$. These are called *triplets* and are denoted respectively by $uv|w$, $uw|v$ and $wv|u$, where $ij|k$ means that the leaf labelled k is incident to the root.

Let $\mathcal{T} = \{T_1, T_2, \dots, T_k\}$ be a collection of unrooted trees, not necessarily on the same set of taxa. The *display graph* of \mathcal{T} is obtained from the disjoint graph union of all trees in \mathcal{T} by identifying vertices with the same label; see Figure 1(b).

Given an undirected graph $G = (V, E)$, a *bag* is simply a subset of V . A *tree decomposition* of G consists of a tree $T_G = (V(T_G), E(T_G))$ where $V(T_G)$ is a collection of bags such that the following holds: (1) every vertex of V is in at least one bag; (2) for each edge $\{u, v\} \in E$, there exists some bag that contains both u and v ; (3) for each vertex $u \in V$, the bags that contain u induce a connected subtree of T_G . The *width* of a tree decomposition is equal to the cardinality of its largest bag, minus 1. The *treewidth* of a graph G is equal to the minimum width, ranging over all possible tree decompositions of G . A tree with at least one edge has treewidth 1. For a fixed value of k one can determine in linear time whether a graph has treewidth at most k [5].

3 Main results

Unless stated otherwise, we assume that $T_1 = (V_1, E_1)$ and $T_2 = (V_2, E_2)$ are both unrooted binary trees on X . Their display graph is denoted by $D = (V, E)$ and R^D denotes the vertex-edge incidence relation in D . We use adj to denote the vertex-vertex adjacency relation in D . Note that $|V| = 3|X| - 4$ and $|E| = 4|X| - 6$.

3.1 TBR / MAF on unrooted trees

We will start by giving the definitions of a TBR move and of the TBR distance between two unrooted binary trees.

Definition 1 (TBR move). *Given an unrooted binary tree T , a tree bisection and reconnection (TBR) move on T consists of removing an edge of T , say $\{u, v\}$, and then reconnecting the subtrees T_u and T_v as follows: subdividing an edge of T_u with a new vertex p ; subdividing an edge of T_v with a new vertex q ; connecting p to q ; and finally suppressing any vertices of degree 2.*

TBR distance is then defined naturally as follows:

Problem: $d_{TBR}(T_1, T_2)$

Input: Two unrooted binary trees T_1, T_2 on the same set of taxa X .

Output: The minimum number of TBR moves required to transform T_1 into T_2 .

We will now give the definition of an uMAF for two unrooted binary trees T_1, T_2 on X . Any collection of trees whose label sets partition X is said to be a *forest on X* . Furthermore, we say that a set $\mathcal{F} = \{F_1, \dots, F_k\}$ of unrooted binary phylogenetic trees – with $|\mathcal{F}|$ referred to as the *size* of \mathcal{F} – is a *forest for T* if \mathcal{F} can be obtained from T by deleting a $(k - 1)$ -sized subset E of $E(T)$, suppressing any unlabeled leaves, and then finally suppressing any vertices with degree 2. To ease reading, we write $\mathcal{F} = T - E$ if \mathcal{F} can be obtained in this way.

Definition 2 (uMAF). *A set \mathcal{F} of unrooted trees is an agreement forest for T_1 and T_2 (denoted *uAF*) if \mathcal{F} is a forest of both T_1 and T_2 . An unrooted maximum agreement forest (uMAF), is an uAF of minimum size.*

So, the uMAF problem is defined as follows:

Problem: $uMAF(T_1, T_2)$

Input: Two unrooted binary trees T_1, T_2 on the same set of taxa X .

Output: An uMAF for T_1 and T_2 .

The two problems defined above are closely related, and known to be NP-hard [1].

Theorem 1 ([1]). *Given two unrooted binary trees T_1, T_2 on the same set of taxa X , we have that $d_{TBR}(T_1, T_2) = |uMAF(T_1, T_2)| - 1$.*

Fortunately, they have been proved to be FPT in their natural parameters [1], and fast algorithms have been recently proposed [27,13]. In this section, we will give a more compact proof of their fixed parameter tractability.

Theorem 2. *Let T_1, T_2 be two unrooted binary trees on the same set of taxa X such that a uAF of size k for these two trees exists. Then, the treewidth of their display graph D is at most $k + 1$.*

Proof. From [19], we know that the display graph of two identical trees has treewidth 2 (or 1 in the case that both trees consist of a single vertex). Thus, if we have an uAF $\mathcal{F} = \{F_1, \dots, F_k\}$ of size k , this means that the display graph D_0 of \mathcal{F} (which we define as the display graph constructed from two disjoint copies of \mathcal{F}) has k connected components, and treewidth at most 2. This is because the treewidth of a disconnected graph is equal to the largest treewidth ranging over its connected components. Now, we can construct a tree decomposition of D from the tree decomposition of \mathcal{F} as follows: suppose \mathcal{F} can be obtained by removing from T_1 , respectively T_2 , a subset of edges K_1 , respectively K_2 , and suppressing vertices with degree 2 and unlabeled leaves. First, note that we can reintroduce the suppressed vertices (and their corresponding edges) in \mathcal{F} , obtaining a new forest \mathcal{F}' , without changing the treewidth. Indeed, given an edge $\{u, v\}$ in \mathcal{F} that corresponded to a path (u, x_1, \dots, x_j, v) before the suppression of the vertices with degree 2, we know that there exists a bag B in the tree decomposition of D_0 such that u and v are in B . Then we can add a set of bags $\{B_1, \dots, B_j\}$ such that $B_1 = \{u, x_1, v\}$, $B_2 = \{x_1, x_2, v\}$, \dots , $B_j = \{x_{j-1}, x_j, v\}$, and add edges $\{B, B_1\}, \{B_1, B_2\}, \dots, \{B_{j-1}, B_j\}$ to the tree decomposition. For the suppressed unlabeled leaves, say u , this is even easier: we add a bag $\{u, v\}$ as child of any of the bags containing v , where v is the vertex from which the suppressed leaf was hanging. It is easy to see that this is a tree decomposition of the display graph of \mathcal{F}' with treewidth 2. Now, we can easily reintroduce the $k - 1$ edges in K_1 to the display graph, again without changing the treewidth, by, for each edge $\{u, v\}$ in K_1 , adding a bag $\{u, v\}$ between two existing bags, one containing u and the other containing v . Note that the obtained decomposition is still a tree, since we are connecting two components of \mathcal{F}' . Now, when adding back the edges of K_2 , this is not true anymore. In this case, there exists at least a path in the tree decomposition, connecting a bag containing u to a bag containing v . Then, taking the shortest of these paths and adding u to its bags not containing u , we increase the treewidth by at most 1. If we do this for all edges in K_2 , we obtain a tree decomposition for the display graph of T_1 and T_2 with treewidth at most $2 + (k - 1) = k + 1$. Note that this bound is tight, as the following example shows: an uMAF of two quartets with different topologies, $uv|wx$ and $ux|vw$ say, contains 2 components, and the display graph of these two quartets has treewidth 3 (see also [19]). \square

In the following, we will demonstrate that $|uMAF(T_1, T_2)| =: k$ can be computed in time $O(f(k) \cdot |X|)$ for some computable function f that depends only on k . We do this via the machinery of MSOL. The high-level idea is that we formulate a logical query to answer the question “Is $k \leq k'?$ ” for increasing values of k' until the answer is *yes*, and then stop: at this point $k' = k$. We use the stronger variant of MSOL that allows quantification over both edges and vertices. In particular, we will use the *extended* MSOL framework of Arnborg et al [2]. Following [10,25] we note that the sets V_1, E_1, V_2, E_2, X (and later, ρ) are all available to the MSOL query i.e. within the query we can distinguish which vertices/edges of D belong to T_1 , which belong to T_2 , and which are taxa.

More formally, we construct an MSOL formula $\Phi(K_1, K_2)$ and a relational structure \mathbf{G} such that $\mathbf{G} \models \Phi(K_1, K_2)$ if and only if K_1 is a set of $k' - 1$ edges of E_1 , and K_2 is a set of $k' - 1$ edges of E_2 , such that, after deleting them, the resulting components form an uAF \mathcal{F} for both T_1 and T_2 , i.e. $\mathcal{F}_1 = T_1 - K_1 = \mathcal{F}_2 = T_2 - K_2$. To model this, we need to have that: (1) the two forests \mathcal{F}_1 and \mathcal{F}_2 induce an identical partition of X and (2) the components of the two induced forests must have the same topology. To enforce (1) we observe that (in, say, T_1) two taxa x_1 and x_2 are in the same component of the forest resulting from deletion of K_1 if and only if they can still reach each other inside T_1 after deletion of those edges. In turn, this occurs if and only if there is a path from x_1 to x_2 entirely contained inside T_1 which avoids all the edges in K_1 . To enforce (2) we demand that a quartet is in the first forest (i.e. the quartet is contained inside one of the trees in the forest) if and only if the quartet is in the second forest. This uses the fact that two unrooted binary trees on the same set of taxa are topologically identical if and only if they induce identical sets of quartets [12].

Before defining $\Phi(K_1, K_2)$, we need to introduce several intermediate predicates. These build on a number of basic predicates which we mainly list for the benefit of readers not familiar with MSOL. They are used to:

- test that Z is equal to the union of two sets P and Q :
 $P \cup Q = Z := \forall z(z \in Z \Rightarrow z \in P \vee z \in Q) \wedge \forall z(z \in P \Rightarrow z \in Z) \wedge \forall z(z \in Q \Rightarrow z \in Z)$
- test that $P \cap Q = \emptyset$:
 $NoIntersect(P, Q) := \forall u \in P(u \notin Q)$
- test that $P \cap Q = \{v\}$:
 $Intersect(P, Q, v) := (v \in P) \wedge (v \in Q) \wedge \forall u \in P(u \in Q \Rightarrow (u = v))$
- test if the sets P and Q are a bipartition of Z :
 $Bipartition(Z, P, Q) := (P \cup Q = Z) \wedge NoIntersect(P, Q)$
- test if the elements in $\{x_1, x_2, x_3, x_4\}$ are pairwise different:
 $allDiff(x_1, x_2, x_3, x_4) := \bigwedge_{i \neq j \in \{1, 2, 3, 4\}} x_i \neq x_j$
- check if the nodes p and q are adjacent in D :
 $adj(p, q) := \exists e \in E(R^D(e, p) \wedge R^D(e, q))$

The predicate $PAC(Z, x_1, x_2, K_i)$ (“*path avoids cuts?*”) asks: is there a path from x_1 to x_2 entirely contained inside vertices Z that avoids all the edges K_i ? We model this by observing that this does *not* hold if you can partition Z into two pieces P and Q , with $x_1 \in P$ and $x_2 \in Q$, such that the only edges that cross the induced cut (if any) are in K_i :

$$\begin{aligned} PAC(Z, x_1, x_2, K_i) := & (x_1 = x_2) \vee \neg \exists P, Q (Bipartition(Z, P, Q) \wedge x_1 \in P \wedge x_2 \in Q \wedge \\ & (\forall p, q(p \in P \wedge q \in Q \Rightarrow \neg adj(p, q) \vee (\exists g \in K_i (R^D(g, p) \\ & \wedge R^D(g, q))))) \end{aligned}$$

We model that a quartet is in the forest (of, say, T_1) by stipulating that there is an embedding (i.e. subdivision) of the quartet, completely contained inside T_1 , which avoids all the edges in K_1 . To model the embedding, we model the five edges of the quartet as five subsets of vertices A, B, C, D, P , representing the subdivisions of the five edges, with P being the central edge and u and v being its endpoints. We demand that (with the exception of u and v) these subsets are disjoint. This is all combined in the following QAC^1 predicate (“*quartet avoids cuts in T_1 ?*”), which returns true if and only if T_1 contains an embedding of $x_a x_b | x_c x_d$ that is disjoint from the

edge cuts K_1 .

$$\begin{aligned}
QAC^1(x_a, x_b, x_c, x_d, K_1) := & \exists u, v \in V_1((u \neq v) \wedge \exists A, B, C, D, P \subseteq V_1(x_a, u \in A \wedge x_b, u \in B \wedge x_c, \\
& v \in C \wedge x_d, v \in D \wedge u \in P \wedge v \in P \wedge Intersect(A, B, u) \wedge \\
& Intersect(A, P, u) \wedge Intersect(B, P, u) \wedge Intersect(C, D, v) \wedge \\
& Intersect(C, P, v) \wedge Intersect(D, P, v) \wedge NoIntersect(A, C) \wedge \\
& NoIntersect(B, C) \wedge NoIntersect(A, D) \wedge NoIntersect(B, D) \wedge \\
& PAC(A, u, x_a, K_1) \wedge PAC(B, u, x_b, K_1) \wedge PAC(C, v, x_c, K_1) \wedge \\
& PAC(D, v, x_d, K_1) \wedge PAC(P, u, v, K_1)))
\end{aligned}$$

We can define $QAC^2(x_a, x_b, x_c, x_d, K_2)$ in a similar way. Note that, for every four taxa, we need to consider all three possible quartet topologies. Then we define $\Phi(K_1, K_2)$ as follows:

$$\begin{aligned}
& \left(\bigwedge_{i \in \{1,2\}} |K_i| = k' - 1 \right) \wedge \left(\bigwedge_{i \in \{1,2\}} K_i \subseteq E_i \right) \wedge \forall x_1, x_2 \in X (PAC(V_1, x_1, x_2, K_1) \Leftrightarrow \\
& PAC(V_2, x_1, x_2, K_2)) \wedge \forall x_1, x_2, x_3, x_4 \in X (allDiff(x_1, x_2, x_3, x_4) \Rightarrow \\
& ((QAC^1(x_1, x_2, x_3, x_4, K_1) \Leftrightarrow QAC^2(x_1, x_2, x_3, x_4, K_2)) \wedge \\
& (QAC^1(x_1, x_3, x_2, x_4, K_1) \Leftrightarrow QAC^2(x_1, x_3, x_2, x_4, K_2)) \wedge \\
& (QAC^1(x_1, x_4, x_2, x_3, K_1) \Leftrightarrow QAC^2(x_1, x_4, x_2, x_3, K_2))))).
\end{aligned}$$

(The cardinality operator is permitted because the extended MSOL framework of [2] allows the incorporation of an *evaluation relation* which can test, amongst other things, the cardinalities of free set variables).

Theorem 3. *Computation of TBR / uMAF on two unrooted binary trees on the same set of taxa X is linear time FPT. That is, the optimum k can be computed in time $O(f(k) \cdot |X|)$ for some computable function that only depends on k .*

Proof. We have presented a logical query to answer the question “Is $k \leq k'$?” for increasing values of k' . For each value of k' the MSOL query, which examines the display graph D , has fixed length. Combining this with the fact that the treewidth of D is bounded by a function of k (by Theorem 2), and that the size of D is a linear function of $|X|$, we have the desired result. (Note that the actual edge cuts - which can be used to construct a uMAF - can also be obtained in the same time bound by leveraging Theorem 4 of [10].) \square

3.2 rSPR / MAF on rooted trees

In this section, we will give a compact proof that the computation of rSPR distance is FPT in its natural parameter. Before that, we need to introduce some definitions.

Definition 3 (rSPR move). *Given a rooted binary tree T , a subtree prune and regraft (rSPR) move on T consists of removing an edge of T , say (u, v) , yielding two trees T_u and T_v , and then reconnecting them as follows: subdividing some edge of T_u with a new vertex p ; adding an edge directed from p to v , and then suppressing any vertices with indegree and outdegree both equal to 1.*

rSPR distance is defined analogously to TBR distance, and a MAF for two rooted binary trees T_1, T_2 is defined similarly to a uMAF, but in a rooted framework. We refer to [6] for precise definitions. The main difference is that a forest consists of *rooted* binary trees and this has to be taken into account when comparing the topology of the components. In the rooted context MAFs are mainly studied because of their close relationship to rSPR distance. To accurately model rSPR distance it is necessary to slightly modify each input tree T_i as follows: we add a vertex with special

label ρ at the end of a pendant edge adjoined to the original root of T_i , see Figure 1(c). We then consider ρ to be part of the label set of the tree. Note that the addition of ρ means that we can equivalently view each T_i as an unrooted binary tree, with ρ acting as a placeholder for the root location, and this is how the trees will be modelled in the display graph.

The close relationship between MAF (assuming ρ has been added as described) and rSPR distance is summarized by the following well-known result.

Theorem 4 ([6]). *Given two rooted binary trees T_1, T_2 on the same set of taxa X , we have that $d_{rSPR}(T_1, T_2) = |MAF(T_1, T_2)| - 1$.*

Note that these problems have been proved NP-hard and FPT in their natural parameter [6].

The MSOL formulation is similar to the TBR formulation, but with the following changes. When checking that the components induced by the edge cuts partition the taxa in the same way in both T_1 and T_2 (i.e. by considering pairs of taxa that still have a path between them), we need to range over $X \cup \{\rho\}$ instead of just X . More significantly, we need predicates for triplets instead of quartets, because we are working in the rooted environment and two rooted binary trees are topologically equivalent if and only if they contain the same set of triplets [11]. Fortunately we can use the fact that triplet $xy|z$ is in T_i ($x, y, z \in X$) if and only if quartet $xy|\rho z$ is in the unrooted interpretation of T_i .

However we cannot simply use ρ as the fourth parameter x_d to QAC^i because this will evaluate to false if the path from ρ to the rest of the quartet embedding has been cut. This is not what we need: ρ is in this context only there to indicate direction, so its particular arm of the quartet embedding can be cut without consequence. We can remedy this by introducing predicates $Quartet^i$ and $Triplet^i$ which check whether the corresponding quartet/triplet was in the *original* tree (i.e. before the edge cuts). We can then leverage the fact that, if three distinct taxa x, y, z are in the same component of the forest, the unique triplet topology they induce within the component will be the same topology as they induced in the original tree.

We first need the following predicate, which is a specialization of the earlier PAC predicate. It tests whether there is a path from x_1 to x_2 that is entirely contained inside vertex set Z :

$$\begin{aligned} path(Z, x_1, x_2) := & (x_1 = x_2) \vee \neg \exists P, Q (Bipartition(Z, P, Q) \wedge x_1 \in P \wedge x_2 \in Q \\ & \wedge (\forall p, q (p \in P \wedge q \in Q \Rightarrow \neg adj(p, q)))) \end{aligned}$$

For each tree T_i , the following predicate checks whether the quartet $x_a x_b | x_c x_d$ is contained in T_i :

$$\begin{aligned} Quartet^i(x_a, x_b, x_c, x_d) := & \exists u, v \in V_i ((u \neq v) \wedge \exists A, B, C, D, P \subseteq V_i (x_a, u \in A \wedge \\ & x_b, u \in B \wedge x_c, v \in C \wedge x_d, v \in D \wedge u \in P \wedge v \in P \wedge Intersect(A, B, u) \wedge \\ & Intersect(A, P, u) \wedge Intersect(B, P, u) \wedge Intersect(C, D, v) \wedge \\ & Intersect(C, P, v) \wedge Intersect(D, P, v) \wedge NoIntersect(A, C) \wedge \\ & NoIntersect(B, C) \wedge NoIntersect(A, D) \wedge NoIntersect(B, D) \wedge \\ & path(A, u, x_a) \wedge path(B, u, x_b) \wedge path(C, v, x_c) \wedge path(D, v, x_d) \wedge \\ & path(P, u, v))) \end{aligned}$$

For each rooted tree T_i , the following predicate checks whether the triplet $x_a x_b | x_c$ is contained in T_i (simply by checking whether $x_a x_b | x_c \rho$ is contained in it):

$$Triplet^i(x_a, x_b, x_c) := Quartet^i(x_a, x_b, x_c, \rho)$$

Now, we are ready to define TAC^i (“*triplet avoids cuts in T_i* ?”), which models whether a triplet is in the forest of T_i induced by the edge cuts:

$$\begin{aligned}
TAC^i(x_a, x_b, x_c, K_1) := & Triplet^i(x_a, x_b, x_c) \wedge \exists u, v \in V_i((u \neq v) \wedge \exists A, B, C, D, P \subseteq V_i(x_a, u \in A \wedge \\
& x_b, u \in B \wedge x_c, v \in C \wedge \rho, v \in D \wedge u \in P \wedge v \in P \wedge Intersect(A, B, u) \wedge \\
& Intersect(A, P, u) \wedge Intersect(B, P, u) \wedge Intersect(C, D, v) \wedge \\
& Intersect(C, P, v) \wedge Intersect(D, P, v) \wedge NoIntersect(A, C) \wedge \\
& NoIntersect(B, C) \wedge NoIntersect(A, D) \wedge NoIntersect(B, D) \\
& \wedge PAC(A, u, x_a, K_1) \wedge PAC(B, u, x_b, K_1) \wedge PAC(C, v, x_c, K_1) \wedge \\
& path(D, v, \rho, K_1) \wedge PAC(P, u, v, K_1)))
\end{aligned}$$

Note how we use *path* rather than *PAC* to model the path from v to ρ i.e. because it does not matter for the triplet whether this path is cut. The final MSOL formulation is then very similar to that given in Section 3.1:

$$\begin{aligned}
(\bigwedge_{i \in \{1,2\}} |K_i| = k' - 1) \wedge (\bigwedge_{i \in \{1,2\}} K_i \subseteq E_i) \wedge \forall x_1, x_2 \in X \cup \{\rho\} (PAC(V_1, x_1, x_2, K_1) \Leftrightarrow \\
PAC(V_2, x_1, x_2, K_2)) \wedge \forall x_1, x_2, x_3 \in X (allDiff(x_1, x_2, x_3) \Rightarrow \\
(TAC^1(x_1, x_2, x_3, K_1) \Leftrightarrow TAC^2(x_1, x_2, x_3, K_2)) \wedge \\
(TAC^1(x_1, x_3, x_2, K_1) \Leftrightarrow TAC^2(x_1, x_3, x_2, K_2)) \wedge \\
(TAC^1(x_2, x_3, x_1, K_1) \Leftrightarrow TAC^2(x_2, x_3, x_1, K_2)))
\end{aligned}$$

Theorem 5. *Computation of rSPR / MAF on two rooted binary trees on the same set of taxa X is linear time FPT. That is, the optimum k can be computed in time $O(f(k) \cdot |X|)$, for some computable function that only depends on k .*

Proof. An agreement forest of the two rooted trees T_1 and T_2 induces an agreement forest (consisting of unrooted binary trees) of the same size of the unrooted interpretations of these trees, simply by ignoring the orientation of edges. Hence the treewidth bound described in Theorem 2 is still applicable, and the theorem follows. (Again, if required one can obtain the actual edge cuts, which can be used to build a MAF, in the same time bound by leveraging Theorem 4 of [10]). \square

3.3 Hybridization Number

In this section, we deal again with rooted trees, and thus we add a vertex labeled ρ to both trees to indicate the root location, as done for rSPR; see Figure 1(c). A *rooted phylogenetic network* (rooted network for short) $N = (V(N), E(N))$ on a set of taxa X is any rooted acyclic digraph in which no vertex has degree 2 (except possibly the root) and whose leaves are bijectively labeled by elements of X . The *hybridization number* of N , denoted by $h(N)$, is defined as

$$h(N) = \sum_{v \in V(N): \delta^-(v) > 0} (\delta^-(v) - 1) = |E(N)| - |V(N)| + 1$$

where $\delta^-(v)$ denotes the indegree of v .

Given a rooted network N on X and a rooted binary tree T on X' , with $X' \subseteq X$, we say that T is *displayed* by N if T can be obtained from N by deleting a subset of its edges and any resulting degree 0 vertices, and then suppressing vertices with $\delta^-(v) = \delta^+(v) = 1$.

We are now ready to define the hybridization number problem:

Problem: $HN(T_1, T_2)$

Input: Two rooted binary trees T_1, T_2 on the same set of taxa X .

Output: A rooted network N displaying T_1 and T_2 such that $h(N)$ is minimum over all rooted networks with this property.

The hybridization number for T_1 and T_2 , denoted by $h(T_1, T_2)$, is defined as the hybridization number of this minimum network. As done for TBR and rSPR, we can give a characterization of the hybridization number in terms of agreement forests. To do so, we need to define *acyclic* agreement forests.

Let $\mathcal{F} = \{F_1, F_2, \dots, F_k\}$ be an agreement forest for two rooted binary trees T_1 and T_2 on the same set of taxa X , and let $AG(T_1, T_2, \mathcal{F})$ be the directed graph whose vertex set is \mathcal{F} and for which (F_i, F_j) is an arc iff $i \neq j$, and either

- (1) the root of $T_1(\mathcal{L}(F_i))$ is an ancestor of the root of $T_1(\mathcal{L}(F_j))$ in T_1 , or
- (2) the root of $T_2(\mathcal{L}(F_i))$ is an ancestor of the root of $T_2(\mathcal{L}(F_j))$ in T_2 .

We call \mathcal{F} an *acyclic agreement forest* (AAF) for T_1 and T_2 if $AG(T_1, T_2, \mathcal{F})$ does not contain any directed cycle. A *maximum* acyclic agreement forest (MAAF), is an AAF of minimum size.

The acyclicity condition is used to model the fact that species cannot inherit genetic material from their own offspring. The two problems defined above are closely related, as the following well-known result shows.

Theorem 6 ([3]). *Given two rooted binary trees T_1, T_2 on the same set of taxa X , we have that $h(T_1, T_2) = |MAAF(T_1, T_2)| - 1$.*

The above equivalence formed the basis for results proving that both problems are NP-hard [8] and fixed parameter tractable [7].

Here we show an alternative proof that computation of hybridization number on two rooted binary trees with the same set of taxa X is FPT, again using MSOL. We will do this by demonstrating that $|MAAF(T_1, T_2)| =: k$ can be computed in time $O(f(k) \cdot |X|)$ for some computable function f that depends only on k . Again, we will formulate a logical query on the display graph to answer the question “Is $k \leq k'$?” for increasing values of k' , until $k' = k$ is reached and the answer to the query is “yes”. Unlike the formulations given earlier for TBR and rSPR, the query has no free variables, and the length of each query will grow as a function of k' . However, given that $k' \leq k$, the length will remain bounded by a function of k . Note that, if a MAAF of size k exists for T_1 and T_2 , then an AF of size k exists too, and as argued for rSPR, if two rooted trees have an agreement forest of size k then so do the underlying, unrooted trees. So the treewidth bound of Theorem 2 is still valid, where $k = |MAAF(T_1, T_2)|$, and this implies that an overall running time of the form $O(f(k) \cdot |X|)$ can be achieved.

The major challenge when modelling MAAF is to encode the acyclicity constraints. It is not clear whether the formulations from the previous sections, in which agreement forests are modelled directly as sets of edge-cuts, can be (elegantly) extended to include acyclicity constraints. For this reason we choose to discard the agreement forest abstraction, using it only to generate the treewidth upper bound. For the actual modelling we use an alternative “elimination-ordering” characterization of MAAF/HN, first presented in [23], which we briefly summarize here.

Given a rooted binary tree T on X , we say a subtree T' of T is *pendant* if there exists a vertex u of T such that $T' = T_u$. In this case it is then natural to associate T' with the subset of X labeling its leaves, i.e. $\mathcal{L}(T')$. We say that T' is a *common pendant* subtree of T_1 and T_2 if it is a pendant subtree of both T_1 and T_2 . We call (S_1, S_2, \dots, S_p) ($p \geq 0$) a *common pendant subtree sequence* of T_1 and T_2 of length p if for every $1 \leq i \leq p$, S_i is a common pendant subtree of $T_1 - \cup_{j < i} \mathcal{L}(S_j)$ and $T_2 - \cup_{j < i} \mathcal{L}(S_j)$. We say that such a sequence is additionally a *tree sequence* if the two trees $T_1 - \cup_{j \leq p} \mathcal{L}(S_j)$ and $T_2 - \cup_{j \leq p} \mathcal{L}(S_j)$ are identical. Informally, a tree sequence of length p describes a sequence of p common pendant subtrees that can be successively pruned from the original trees to reach a common core tree. If T_1 and T_2 are already identical then we use the empty tree sequence \emptyset , and take $p = 0$, to represent this.

The results in [23] establish that $h(T_1, T_2)$ is equal to the smallest p such that a tree sequence of length p exists. This is the characterization of optimality that we will use i.e. each logical query will

pose the question, “Does a tree sequence of length k' exist?”. There is no need to model acyclicity in this formulation. However, we do need to model the concept *common pendant subtree* and the impact of earlier pruning steps on the original trees.

Before writing down the MSOL formulation we need some new auxiliary predicates. The first predicate checks whether there is a path from x_1 to x_2 within Z that survives the deletion of vertex u . This is similar to the *PAC* predicate defined earlier.

$$\begin{aligned} \text{pathSurvivesVertexCut}(Z, x_1, x_2, u) := & (u \neq x_1) \wedge (u \neq x_2) \wedge \\ & ((x_1 = x_2) \vee \neg \exists P, Q (Bipartition(Z, P, Q) \wedge \\ & x_1 \in P \wedge x_2 \in Q \wedge (\forall p, q (p \in P \wedge q \in Q \Rightarrow \\ & \neg adj(p, q) \vee p = u \vee q = u)))) \end{aligned}$$

For a vertex $u \neq \rho$ in a tree T_i and a taxon $x \in X$, observe that x is in the clade rooted at u (i.e. in the label set of the pendant subtree rooted at u) if and only if $(x = u)$ or deleting u from T_i destroys all paths from ρ to x (inside T_i). Hence:

$$\text{InCladeUnder}^i(u, x) := (u = x) \vee \neg \text{pathSurvivesVertexCut}(V_i, \rho, x, u)$$

This leads naturally to a predicate for testing whether $C \subseteq X$ is a clade of T_i :

$$\text{Clade}^i(C) := \exists u \in V_i (\forall x (x \in C \Leftrightarrow \text{InCladeUnder}^i(u, x)))$$

As we shall see, it is useful to extend this predicate with an optional list Z_1, Z_2, \dots which represent subsets of X describing common pendant subtrees that have already been pruned from the tree. The statement $\text{Clade}^i(C, Z_1, Z_2, \dots)$ evaluates to true if and only if C is a clade of T_i *after* the taxa in Z_1, Z_2, \dots have been pruned away. (To avoid ambiguity the predicate automatically returns false if C intersects with any of the Z_i .) Note that the list of Z_i is shown in square brackets to emphasize that it is a “macro”: there will be a different predicate for each possible list length t . The list of Z_i will never be longer than $h(T_1, T_2)$, and length of the generated predicate will be bounded by a function of the list length, so the length of the overall logical query remains bounded by a function of $h(T_1, T_2)$.

$$\begin{aligned} \text{Clade}^i(C, [Z_1, \dots, Z_t]) := & \text{NoIntersect}(C, Z_1) \wedge \dots \wedge \text{NoIntersect}(C, Z_t) \wedge \exists u \in V_i \\ & (\forall x (x \in C \Rightarrow \text{InCladeUnder}^i(u, x))) \wedge \\ & \forall x (\text{InCladeUnder}^i(u, x) \Rightarrow (x \in C) \vee (x \in Z_1) \vee \dots \vee (x \in Z_t)) \end{aligned}$$

We are now ready to define the CPS (i.e. “common pendant subtree”) predicate. We do this by observing that $C \subseteq X$ corresponds to a common pendant subtree of T_1 and T_2 if and only if C is a clade of both trees (this ensures that C is pendant in both trees) and the set of triplets induced by C is identical in both trees (this ensures that the pendant subtree has the same topology in both trees).

$$\begin{aligned} \text{CPS}(T_1, T_2, C) := & \text{Clade}^1(C) \wedge \text{Clade}^2(C) \wedge \forall x \forall y \forall z (x, y, z \in C \wedge \\ & \text{allDiff}(x, y, z) \Rightarrow (\text{Triplet}^1(x, y, z) \Leftrightarrow \text{Triplet}^2(x, y, z))) \end{aligned}$$

We extend this now with a list of Z_i representing the taxa we have already pruned. This new version of the predicate evaluates to true if and only if C corresponds to a common pendant subtree in the two trees *after* all the Z_i have been pruned away. (Here we make implicit use of the fact that the *Clade* predicate immediately returns false whenever C intersects with the Z_i .)

$$\begin{aligned} \text{CPS}(T_1, T_2, C, [Z_1, \dots, Z_t]) := & \text{Clade}^1(C, Z_1, \dots, Z_t) \wedge \text{Clade}^2(C, Z_1, \dots, Z_t) \wedge \forall x \forall y \forall z \\ & (x, y, z \in C \wedge \text{allDiff}(x, y, z) \Rightarrow (\text{Triplet}^1(x, y, z) \Leftrightarrow \\ & \text{Triplet}^2(x, y, z))) \end{aligned}$$

We are now ready to directly pose the question: is there a tree sequence of length k' ? We can assume $k' \geq 1$ because $k' = 0$ is trivial to check in polynomial time. To make the formulation slightly more compact we actually construct a list of length $k' + 1$, where $C_{k'+1}$ represents the taxa that still remain after the common pendant subtrees have been pruned away: we can then test that the sequence is a *tree* sequence (i.e. that a common core tree remains) by testing that $CPS(T_1, T_2, C_{k'+1}, C_1, \dots, C_{k'})$ is true. Note that the *HybNum* predicate is again a macro, whose expansion depends on k' .

$$\begin{aligned} HybNum[k'](T_1, T_2) := & \exists C_1, \dots, C_{k'}, C_{k'+1} (Partition(X, C_1, \dots, C_{k'+1}) \wedge \\ & CPS(T_1, T_2, C_1, \emptyset) \wedge \\ & CPS(T_1, T_2, C_2, C_1) \wedge \\ & CPS(T_1, T_2, C_3, C_1, C_2) \wedge \\ & \dots \\ & CPS(T_1, T_2, C_{k'+1}, C_1, \dots, C_{k'})). \end{aligned}$$

The *Partition* predicate has the expected meaning and definition:

$$Partition(X, C_1, \dots, C_k) := (\wedge_{i \neq j} NoIntersect(C_i, C_j)) \wedge \forall u (u \in X \Leftrightarrow (u \in C_1 \vee u \in C_2 \vee \dots \vee u \in C_k))$$

Concluding, we have the following result:

Theorem 7. *Computation of hybridization number / MAAF on two rooted binary trees on the same set of taxa X is linear time FPT. That is, the optimum k can be computed in time $O(f(k) \cdot |X|)$, for some computable function that only depends on k .*

3.4 Parsimony distance on binary characters

Let T be an unrooted binary tree on a set of taxa X . A *binary character* f is simply a function $f : X \rightarrow \{red, blue\}$. An *extension* of f to T is a mapping $g : V(T) \rightarrow \{red, blue\}$ such that, for all $x \in X$, $g(x) = f(x)$. For a given character f , an *optimal extension* is any extension g of f such that the number of bichromatic edges is minimized. The number of bichromatic edges in an optimal extension is called the *parsimony score* of f with respect to T , and denoted $l_f(T)$. The well-known algorithm by Fitch can be used to compute $l_f(T)$ (and an optimal extension) in polynomial time [18]. We shall describe Fitch's algorithm in due course. The *parsimony distance problem on binary characters*, denoted d_{MP}^2 , is defined as follows [17].

Problem: $d_{MP}^2(T_1, T_2)$

Input: Two unrooted binary trees T_1, T_2 on the same set of taxa X

Output: Construct a binary character f on X such that the value $|l_f(T_1) - l_f(T_2)|$ is maximized.

We use d_{MP}^2 to denote the optimum value of $|l_f(T_1) - l_f(T_2)|$. The problem was recently shown to be NP-hard and APX-hard [22]. It is not known whether the problem is FPT in d_{MP}^2 . The following result, however, is already known.

Lemma 1 ([17]). *Let T_1, T_2 be two unrooted binary trees on the same set of taxa X . Then $d_{MP}^2(T_1, T_2) \leq d_{TBR}(T_1, T_2)$.*

Given two trees T_1, T_2 as input to d_{MP}^2 , it is not known whether the display graph D of T_1 and T_2 has treewidth bounded by a function of d_{MP}^2 . However, from Lemma 1 and earlier results in this article (Theorems 1 and 2) it is clear that D has treewidth bounded by a function of $d_{TBR}(T_1, T_2)$. An MSOL formulation modelling d_{MP}^2 , whose length is bounded by a function of d_{MP}^2 , will therefore give a running time of the form $f(d_{TBR}(T_1, T_2)) \cdot O(|X|)$ for some computable function f that only depends on $d_{TBR}(T_1, T_2)$. We now give such a formulation. We will remain

within the framework of [2], this time using the (“linear extremum”) optimization variant of MSOL. This allows us to maximize or minimize an affine function of (the cardinalities of) the free set variables in the query.

The MSOL formulation we give here, which is based on an ILP formulation from [22], maximizes $l_f(T_1) - l_f(T_2)$. (To compute d_{MP}^2 we need to use the MSOL machinery twice, once for $l_f(T_1) - l_f(T_2)$ and once for $l_f(T_2) - l_f(T_1)$, taking the maximum of the two results. The second call only differs in its objective function so we omit details).

The basic idea is to range over all possible binary characters, simultaneously embedding two static formulations⁴ of Fitch’s algorithm to “compute” $l_f(T_1) - l_f(T_2)$.

Fitch’s algorithm proceeds as follows. If T is not rooted, we root it arbitrarily (by subdividing an arbitrary edge). The algorithm then works in two phases, a bottom-up phase which computes $l_f(T)$, and then a top-down phase which actually computes a corresponding extension. In the bottom-up phase, we start by assigning each taxon x the singleton set of colours $S(x) := \{f(x)\}$. For an internal node u with children v_1, v_2 we set $S(u) := S(v_1) \cap S(v_2)$ (if $S(v_1) \cap S(v_2) \neq \emptyset$, in which case we say u is an *intersection node*) and $S(u) := S(v_1) \cup S(v_2)$ (if $S(v_1) \cap S(v_2) = \emptyset$, in which case we say that u is a *union node*). The value $l_f(T)$ is equal to the number of internal nodes that are union nodes. (We omit a description of the constructive top-down phase as it is not relevant for this article).

To translate this into an MSOL formulation, we begin by arbitrarily rooting T_1 and T_2 and using ρ as the placeholder for the root, in the usual fashion. The central idea is to partition the vertices of each tree T_i into four possible subsets R^i, B^i, RB_I^i and RB_U^i corresponding to the set of colours that Fitch allocates to each node, and distinguishing union events from intersection events: *red*, *blue*, $\{\text{red}, \text{blue}\}$ (intersection node) and $\{\text{red}, \text{blue}\}$ (union node). We therefore ask the MSOL formulation to instantiate the free set variables R^i, B^i, RB_I^i and RB_U^i ($i \in \{1, 2\}$) such that the expression $|RB_U^1| - |RB_U^2|$ is maximized. (If desired, this can then be made constructive via Theorem 4 of [10].) The only significant work is simulating the bottom-up execution of Fitch’s algorithm. In particular, encoding expressions which describe the state of a parent node u in terms of its two children v_1, v_2 .

We introduce the auxiliary predicate $child^i(u, v)$ which says that v is a child of u in T_i . We can model this as follows: v is a child of u in T_i if and only if there is an edge e in T_i such that v and u are both endpoints of e and there does *not* exist a path from ρ to v that survives the edge cut e . (Here we have specialized the PAC predicate from earlier so that it only takes a single edge, rather than a set of edges, as its fourth argument.)

$$child^i(u, v) := (u \neq v) \wedge \exists e \in E_i(R^D(e, u) \wedge R^D(e, v) \wedge \neg PAC(V_i, \rho, v, e))$$

For each tree T_i we add the following constraints, which encode (in this order):

- The four subsets R, B, RB_I and RB_U partition the vertices of the tree;
- A vertex in X can only be in R or B ;
- An internal node is in R if and only if (one child is in R and the other child is not in B);
- An internal node is in B if and only if (one child is in B and the other child is not in R);
- An internal node is in RB_I if and only if (neither child is in R or B);
- An internal node is in RB_U if and only if (one child is in R and one child is in B).

⁴ Interestingly, the earlier phylogenetics MSOL articles [10,25] also used static formulations: in that case the classical polynomial-time algorithm of Aho.

$$\begin{aligned}
& \left(\text{Partition}(V_i, R^i, B^i, RB_I^i, RB_U^i) \right) \wedge \left(\forall x \in X (x \notin RB_I^i \wedge x \notin RB_U^i) \right) \wedge \\
& \left(\forall u \in V_i (u \notin X \Rightarrow (u \in R^i \Leftrightarrow \exists c_1, c_2 \in V_i ((c_1 \neq c_2) \wedge \text{child}^i(u, c_1) \wedge \text{child}^i(u, c_2) \wedge \right. \\
& \qquad \qquad \qquad \left. c_1 \in R^i \wedge c_2 \notin B^i))) \right) \wedge \\
& \left(\forall u \in V_i (u \notin X \Rightarrow (u \in B^i \Leftrightarrow \exists c_1, c_2 \in V_i ((c_1 \neq c_2) \wedge \text{child}^i(u, c_1) \wedge \text{child}^i(u, c_2) \wedge \right. \\
& \qquad \qquad \qquad \left. c_1 \in B^i \wedge c_2 \notin R^i))) \right) \wedge \\
& \left(\forall u \in V_i (u \notin X \Rightarrow (u \in RB_I^i \Leftrightarrow \exists c_1, c_2 \in V_i ((c_1 \neq c_2) \wedge \text{child}^i(u, c_1) \wedge \text{child}^i(u, c_2) \wedge \right. \\
& \qquad \qquad \qquad \left. c_1 \notin R^i \wedge c_1 \notin B^i \wedge c_2 \notin R^i \wedge c_2 \notin B^i))) \right) \wedge \\
& \left(\forall u \in V_i (u \notin X \Rightarrow (u \in RB_U^i \Leftrightarrow \exists c_1, c_2 \in V_i ((c_1 \neq c_2) \wedge \text{child}^i(u, c_1) \wedge \text{child}^i(u, c_2) \wedge \right. \\
& \qquad \qquad \qquad \left. c_1 \in R^i \wedge c_2 \in B^i))) \right).
\end{aligned}$$

Finally, we ensure that both trees select the same character as follows:

$$\forall x \in X ((x \in R^1 \Leftrightarrow x \in R^2) \wedge (x \in B^1 \Leftrightarrow x \in B^2))$$

This concludes the formulation. Then we have the following result:

Theorem 8. $d_{MP}^2(T_1, T_2)$ is linear time fixed parameter tractable in parameter $d_{TB}^2(T_1, T_2)$.

4 Conclusion

We have demonstrated how agreement forests, which are intensively studied objects in the phylogenetics literature, naturally lead to bounded treewidth in an auxiliary graph structure known as the *display graph*. This opens the door to compact, “declarative” proofs of fixed parameter tractability for a range of phylogenetics problems by formulating them in Monadic Second Order Logic (MSOL). Our formulations have introduced a number of logical predicates and design principles that will hopefully be of use to other phylogenetics researchers seeking to utilize this powerful machinery elsewhere in phylogenetics. Indeed, a natural follow-up question is to ask: what are the essential characteristics of phylogenetics problems that are amenable to this technique?

5 Acknowledgements

We thank Mathias Weller for helpful conversations.

References

1. B.L. Allen and M. Steel. Subtree transfer operations and their induced metrics on evolutionary trees. *Annals of Combinatorics*, 5:1–15, 2001.
2. S. Arnborg, J. Lagergren, and D. Seese. Easy problems for tree-decomposable graphs. *Journal of Algorithms*, 12:308 – 340, 1991.
3. M. Baroni, S. Grünwald, V. Moulton, and C. Semple. Bounding the number of hybridisation events for a consistent evolutionary history. *Mathematical Biology*, 51:171–182, 2005.

4. H. L. Bodlaender. A tourist guide through treewidth. *Acta cybernetica*, 11(1-2):1, 1994.
5. H. L. Bodlaender. A linear-time algorithm for finding tree-decompositions of small treewidth. *SIAM Journal of Computing*, 25:1305–1317, 1996.
6. M. Bordewich and C. Semple. On the computational complexity of the rooted subtree prune and regraft distance. *Annals of Combinatorics*, 8:409–423, 2004.
7. M. Bordewich and C. Semple. Computing the hybridization number of two phylogenetic trees is fixed-parameter tractable. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 4:458–466, 2007.
8. M. Bordewich and C. Semple. Computing the minimum number of hybridization events for a consistent evolutionary history. *Discrete Applied Mathematics*, 155:914–928, April 2007.
9. D. Bryant. The splits in the neighborhood of a tree. *Annals of Combinatorics*, 8(1):1–11, 2004.
10. D. Bryant and J. Lagergren. Compatibility of unrooted phylogenetic trees is FPT. *Theoretical Computer Science*, 351:296 – 302, 2006.
11. D. Bryant and M. Steel. Extension operations on sets of leaf-labeled trees. *Advances in Applied Mathematics*, 16(4):425–453, 1995.
12. O. P. Buneman. The recovery of trees from measures of dissimilarity. *Mathematics in the archaeological and historical sciences*, 1971.
13. J. Chen, J-H. Fan, and S-H. Sze. Parameterized and approximation algorithms for maximum agreement forest in multifurcating trees. *Theoretical Computer Science*, 562:496–512, 2015.
14. B. Courcelle. The monadic second-order logic of graphs. I. Recognizable sets of finite graphs. *Information and Computation*, 85:12–75, 1990.
15. R. G. Downey and M. R. Fellows. *Fundamentals of parameterized complexity*, volume 4. Springer, 2013.
16. J. Felsenstein. *Inferring Phylogenies*. Sinauer Associates, Incorporated, 2004.
17. M. Fischer and S. Kelk. On the maximum parsimony distance between phylogenetic trees. *Annals of Combinatorics*, 2014. preliminary version arXiv preprint arXiv:1402.1553.
18. W. Fitch. Toward defining the course of evolution: minimum change for a specific tree topology. *Systematic Biology*, 20(4):406–416, 1971.
19. A. Grigoriev, S. Kelk, and N. Lekić. On low treewidth graphs and supertrees. In *Algorithms for Computational Biology*, pages 71–82. Springer, 2014.
20. D. H. Huson, R Rupp, and C Scornavacca. *Phylogenetic networks: Concepts, Algorithms and Applications*. Cambridge University Press, 2010.
21. D. H. Huson and C. Scornavacca. A survey of combinatorial methods for phylogenetic networks. *Genome biology and evolution*, 3(1):23–35, January 2011.
22. S. Kelk and M. Fischer. On the complexity of computing mp distance between binary phylogenetic trees. *arXiv preprint arXiv:1412.4076*, 2014.
23. S. Kelk, C. Scornavacca, and L. Van Iersel. On the elusiveness of clusters. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 9(2):517–534, 2012.
24. N. Robertson and P. D. Seymour. Graph minors. II. Algorithmic aspects of tree-width. *Journal of algorithms*, 7(3):309–322, 1986.
25. C. Scornavacca, L. van Iersel, S. Kelk, and D. Bryant. The agreement problem for unrooted phylogenetic trees is FPT. *Journal of Graph Algorithms and Application*, 18:385–392, 2014.
26. C. Semple and M. Steel. *Phylogenetics*. Oxford University Press, 2003.
27. C. Whidden, R. G. Beiko, and N. Zeh. Fixed-parameter algorithms for maximum agreement forests. *SIAM Journal on Computing*, 42(4):1431–1466, 2013.
28. C. Whidden, N. Zeh, and R. G. Beiko. Supertrees based on the subtree prune-and-regraft distance. 63(4):566–581, 2014.